

Computing similarity between text and human summaries. A comparison of models

Marianela Crissman

Motivation

To increase education

- Online Learning
- Two-way communication => Scalability issue
- Intelligent Tutoring Systems



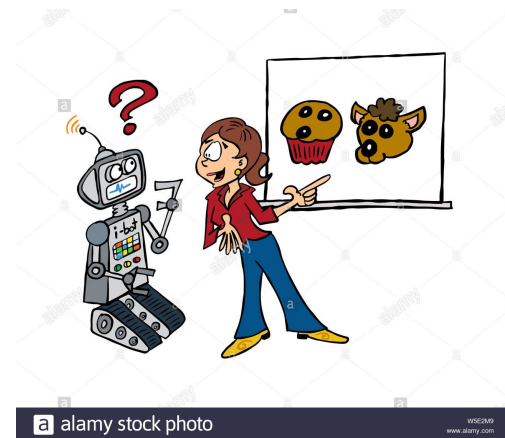
coursera

UDACITY

Khan Academy

YouTube

Comparing user-response vs text book-response.



alamy stock photo

1502161
www.alamy.com

Background

Vector Modeling

1. Corpus

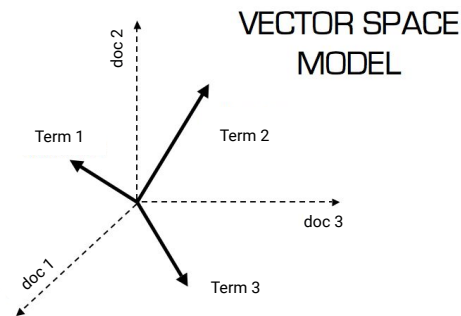
- d1 = "computer master tutor"
d2 = "computer system learning response"
d3 = "learning tutor system"
:

2. Bag of words: unique tokens

3. Term x Document matrix

Similarity computed using Cosine between vectors

High dimension sparse matrices, certain relationships are not captured.



	d1	d2	d3	d4	d5	d6	d7
master	1	0	0	1	0	0	0
tutor	1	0	1	0	0	0	0
computer	1	1	0	0	0	0	0
learning	0	1	1	0	1	0	0
system	0	1	1	2	0	0	0
response	0	1	0	0	1	0	0

$$\text{similarity}(A,B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

Background

Latent Semantic Analysis (LSA)

- Words that appear together are somehow related. LSA captures this relationship and includes it into the vector model.

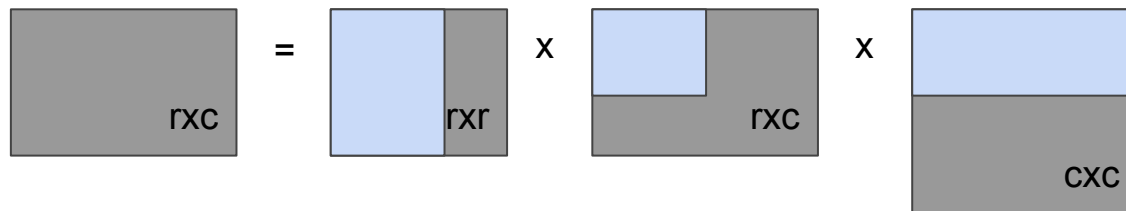
Work on Latent Semantic Analysis (LSA):

- Deerwester et al (1990). Landauer et al (1998) - Latent Semantic Analysis
- Steinberger et al (2012) - Evaluation Measure for Text Summarization
- Nye, B.D., Graesser (2014) - 17 years of AutoTutor

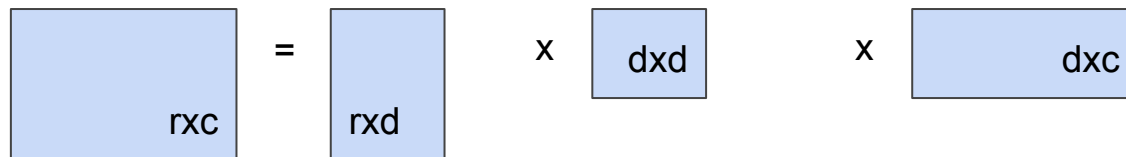
Background

LSA uses dimensionality reduction. Good enough: $d=100..400$.

Singular Value Decomposition (SVD)

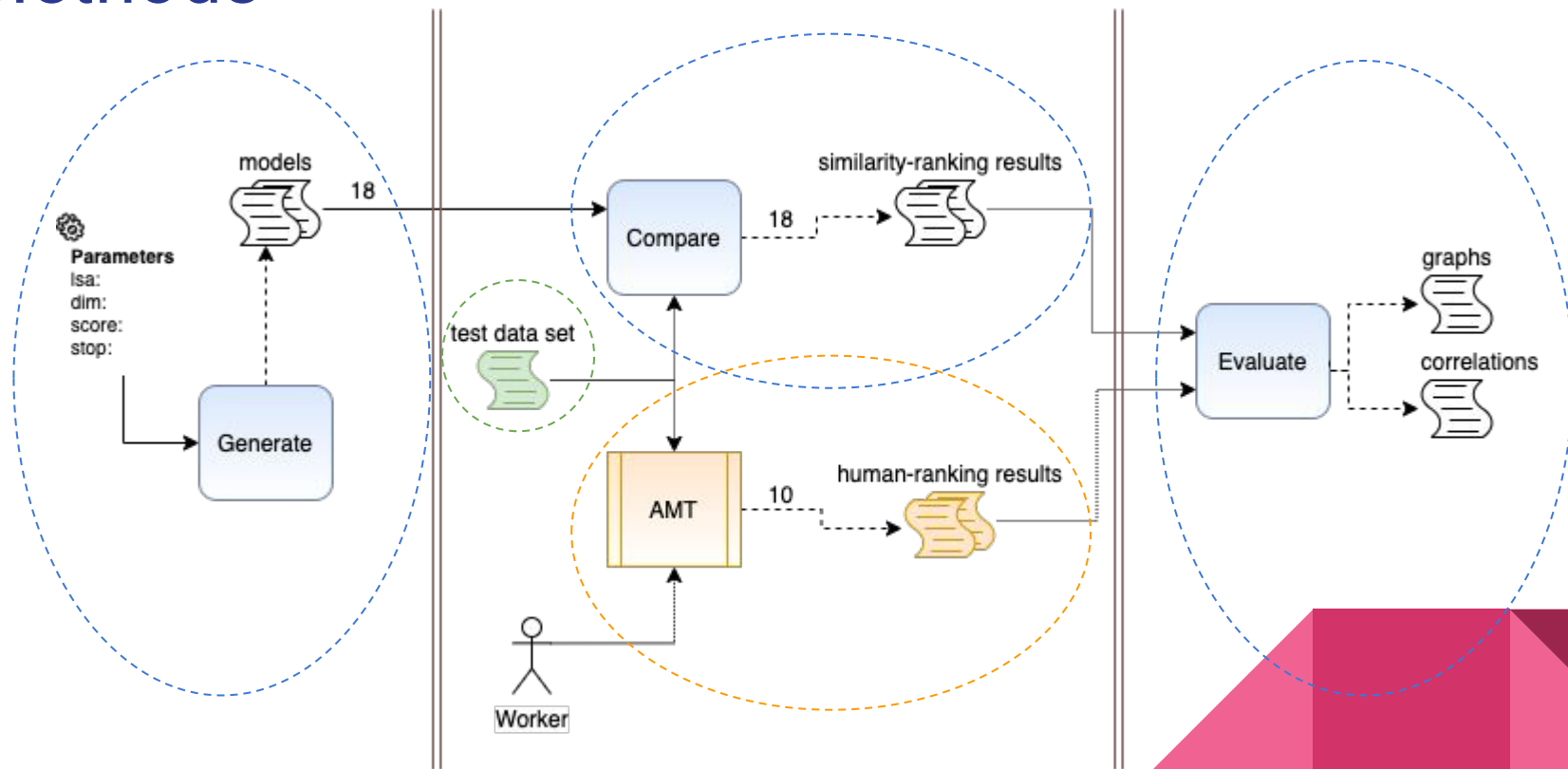


Reconstruction

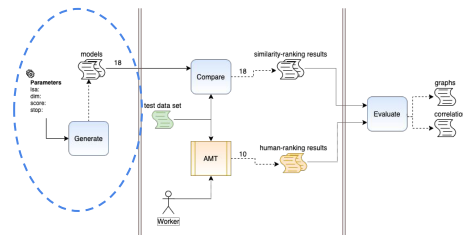


Words that often appeared together in documents would have vectors pointing to similar directions.

Methods



Methods - Generate



1. Corpus: a collection of senate speeches. Data from (Diermeier et al. 2012)
2. Pre-process: Punctuation, Lowercase, @YEAR@ and @PRICE@
3. Vectorization:
 - a. Bag of words: [a-zA-Z0-9]{2,}; unigrams; appeared in at least 2 documents
 - b. Score: frequency, zero-one, TF-IDF
 - c. Stop words: remove words that do not add meaning
- 4?. Apply LSA: d=100 or d=300 dimensions

	d1	d2	d3	d4	d5	d6	d7
master	1	0	0	1	0	0	0
tutor	1	0	1	0	0	0	0
computer	1	1	0	0	0	0	0
learning	0	1	1	0	1	0	0
system	0	1	1	2	0	0	0
response	0	1	0	0	1	0	0

Methods - Generate

VSM	dimensionality	score	stop words
1	-	count	yes
2	-	count	-
3	-	tfidf	yes
4	-	tfidf	-
5	-	zero-one	yes
6	-	zero-one	-
7	100	count	yes
8	100	count	-
9	100	tfidf	yes
10	100	tfidf	-
11	100	zero-one	yes
12	100	zero-one	-
13	300	count	yes
14	300	count	-
15	300	tfidf	yes
16	300	tfidf	-
17	300	zero-one	yes
18	300	zero-one	-

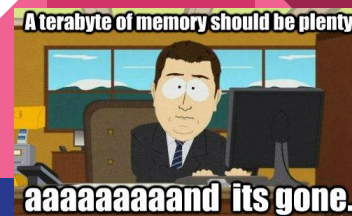
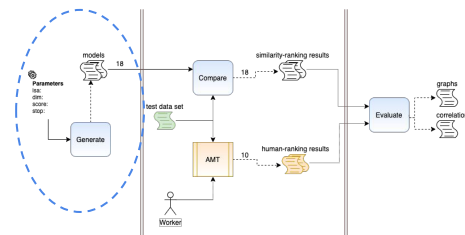
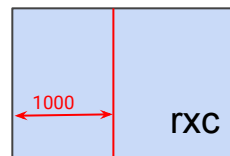
Vector Space Models

- lsa: used lsa or not
- dim: 100 or 300
- sco: count, tfidf, zero-one
- sw: used stop words or not

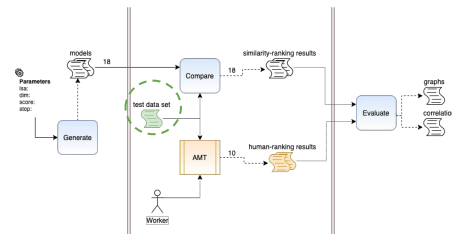
Size of a model: (60k, 130k) ~ 32.5GB

Storage:

- Models 1-6 as sparse matrix.
- Models 7-18, kept first **1000** columns.



Methods - Test Data Set



From *Mi Guia**, we extracted 10 lead paragraphs (LP_x). $x=1..10$
Six volunteers derived summaries (S) for each of the LP_x .

- $S_{y,x}$ = Summary of volunteer y on LP_x

To test

We grouped each LP_x with 6 Summary Paragraphs ($SP(i,x)$, $i=1..6$), where $SP(i,x)$:

- was not similar for $i=1,2$
- was similar for $i=3,4,5,6$

Table 3.2.1. Example of a task to be ranked.

Paragraph: LP1	Cancer begins in cells. Cells are the building blocks that make up all parts of our bodies, including our breasts. Cancer begins when the cells in the breast change and grow into a tumor. If not removed or treated, a breast cancer tumor can spread to other parts of the body and become lethal. These days, there are many treatment options for breast cancer and the majority of women are cured of breast cancer.
sentence SP(1,1)	Systemic therapy reduces the chances of breast cancer coming back.
sentence SP(2,1)	Mastectomy is a procedure that eliminates the breast.
sentence SP(3,1)	Cancer starts in the smallest entity of our bodies; the cell grows too the parts of the body and if left untreated will probably end in death.
sentence SP(4,1)	Cancer is a pathologic function of cells, your body is made by cells. This means that you can have cancer in a specific part of your body or you can have cancer everywhere. When it starts, it is usually in a determined territory and it is easier to treat and get a better prognosis in that moment.
sentence SP(5,1)	Cancer starts in the body cells and when it grows out of control, it becomes a tumor. For breast cancer there are actually several treatment options that can prevent the spreading of malignant cells to other body organs, so a big number of women are cured from this disease.
sentence SP(6,1)	Sometimes cells in the body can change, and may grow into cancerous tumors, which must be removed from the body before they become lethal.

* An application for Breast Cancer education.

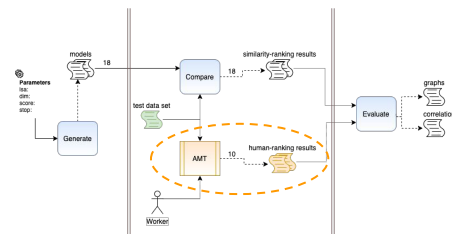
Methods - Human's perception

From each LPx, we created one Human Intelligence Tasks (HIT) in Amazon Mech. Turk

- Each HIT completed by 10 humans

Task: order the blocks depending on how similar they consider these sentences are to the paragraph.

They ended up ranking SP(1,x) - SP(6,x).



Paragraph

For most types of cancer, doctors use staging information to help plan treatment and to predict a person's outlook (prognosis). Although each person's situation is different, cancers with the same stage tend to have similar outlooks and are often treated the same way. The cancer stage is also a way for doctors to describe the extent of the cancer when they talk with each other about a person's cancer.

e1

Doctors usually use stages for cancer so they can predict the evolution and can choose the treatment. Actually they can describe the extension of the tumor through staging.

e2

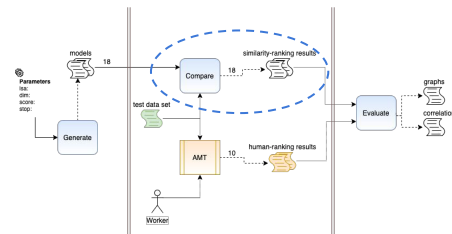
The advancement of cancer is categorized and described in stages, which also works as a guide for treatment and prediction.

e6

When doctors are dealing with cancer patients, Oncologists use staging information as a reference to determine the severity of the patient, and their treatment as well. Even though each patient has their own characteristics, there are similar symptoms during each disease's stage.

Submit

Methods - Automatic Similarity

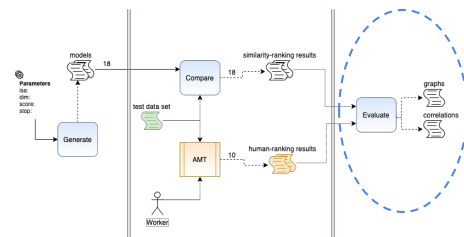


Given a paragraph, a summary paragraph and a model: **LP_x**, **SP(i,x)** and **m**.

1. Preprocessed texts:
 - {?!, lower, @PRICE@, @YEAR@}
2. Vectorized the sentence, using model **m**
 - Sentence vector = mean(word vectors).
3. Compared LP_x paragraph to each one of its SP(i,x),
 - Cosine Similarity(LP_x, SP(i,x)) , for x=1..10 and i=1..6

Comparisons belonging to (**LP_x**, **m**) were sorted and ranked:
1- most similar and 6- least similar.

Methods - Evaluation



Using model $m=1..18$ and paragraph $x=1..10$.

Machine Rankings, $mr_{m,x} = MRanking(m, LPx)$. Human Rankings, $hr_x = HRanking(LPx)$

For each ranking pair ($mr_{m,x}$, hr_x):

- Spearman's Rank coefficient
- Kendall Tau Correlation
- Order Preservation Measure

$$r_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}, \quad d_i = \text{rg}(X_i) - \text{rg}(Y_i)$$

$$\tau = \frac{2}{n(n-1)} \sum_{i < j} \text{sgn}(x_i - x_j) \text{sgn}(y_i - y_j).$$

$$\text{orderPreservationMeasure} = \frac{1}{C_{n,2}} \sum_{i=0}^{n-1} \sum_{j=i+1}^n \text{isAscending}(i, j, A_{\text{num}}),$$

$$\text{isAscending}(i, j, A_{\text{num}}) = 1 \quad \text{if } A_{\text{num}}[i] < A_{\text{num}}[j], \quad 0 \text{ otherwise.}$$

RESULTS

Results - Models

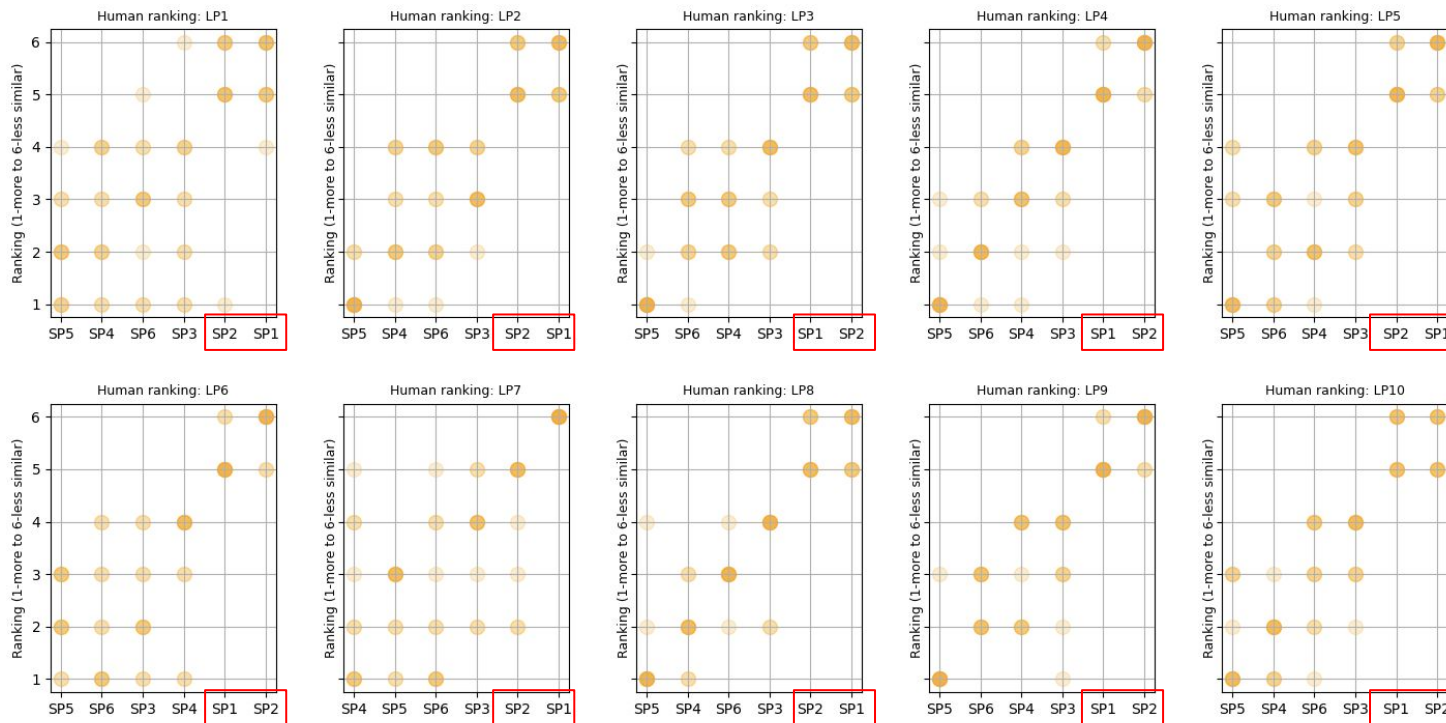
model	
lsa:- dim:- sco:count sw:yes	1
lsa:- dim:- sco:count sw:-	2
lsa:- dim:- sco:tfidf sw:yes	3
lsa:- dim:- sco:tfidf sw:-	4
lsa:- dim:- sco:zeroone sw:yes	5
lsa:- dim:- sco:zeroone sw:-	6
lsa:yes dim:100 sco:count sw:yes	7
lsa:yes dim:100 sco:count sw:-	8
lsa:yes dim:100 sco:tfidf sw:yes	9
lsa:yes dim:100 sco:tfidf sw:-	10
lsa:yes dim:100 sco:zeroone sw:yes	11
lsa:yes dim:100 sco:zeroone sw:-	12
lsa:yes dim:300 sco:count sw:yes	13
lsa:yes dim:300 sco:count sw:-	14
lsa:yes dim:300 sco:tfidf sw:yes	15
lsa:yes dim:300 sco:tfidf sw:-	16
lsa:yes dim:300 sco:zeroone sw:yes	17
lsa:yes dim:300 sco:zeroone sw:-	18

Generated 18 models

Each Vector Space Model:

- In approximately 5 minutes
- Each model took approximately 300 MB on disk
- Each model contained around 60k words.

Results - Human Rankings



Humans ranking of SP(i,x) per LPx.

- Agreed on SP1, SP2 less similar
- Confusion with SP3-SP6.

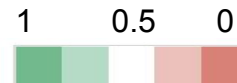
Opacity: 1- more votes, 0- less votes.

1 0.5 0

[illegible]

- Two **not related**.
- Four **related**.

Results - Similarities



Similarities on LP1. SP(i,1) i=1..6.

- Alternation:
 - Stop words remove noise.
- Model's threshold:
 - Model 11: values 0.70 - 0.92
 - Model 3: values 0.30 - 0.65

Each model has a different threshold for determining when something is similar or not.

model		LP1					
		SP5	SP4	SP3	SP6	SP1	SP2
lsa:- dim:- sco:count sw:yes	1	0.804	0.666	0.662	0.620	0.517	0.378
lsa:- dim:- sco:count sw:-	2	0.998	0.996	0.990	0.992	0.980	0.991
lsa:- dim:- sco:tfidf sw:yes	3	0.656	0.450	0.510	0.480	0.419	0.306
lsa:- dim:- sco:tfidf sw:-	4	0.993	0.988	0.975	0.975	0.956	0.975
lsa:- dim:- sco:zeroone sw:yes	5	0.750	0.658	0.695	0.686	0.428	0.253
lsa:- dim:- sco:zeroone sw:-	6	0.992	0.989	0.977	0.977	0.954	0.964
lsa:yes dim:100 sco:count sw:yes	7	0.926	0.868	0.786	0.826	0.615	0.407
lsa:yes dim:100 sco:count sw:-	8	0.998	0.996	0.990	0.991	0.979	0.989
lsa:yes dim:100 sco:tfidf sw:yes	9	0.829	0.830	0.738	0.725	0.525	0.444
lsa:yes dim:100 sco:tfidf sw:-	10	0.996	0.991	0.981	0.984	0.962	0.979
lsa:yes dim:100 sco:zeroone sw:yes	11	0.923	0.924	0.918	0.881	0.824	0.700
lsa:yes dim:100 sco:zeroone sw:-	12	0.997	0.994	0.992	0.986	0.976	0.974
lsa:yes dim:300 sco:count sw:yes	13	0.854	0.642	0.698	0.707	0.490	0.362
lsa:yes dim:300 sco:count sw:-	14	0.998	0.995	0.990	0.991	0.979	0.989
lsa:yes dim:300 sco:tfidf sw:yes	15	0.802	0.800	0.627	0.701	0.473	0.308
lsa:yes dim:300 sco:tfidf sw:-	16	0.995	0.990	0.980	0.981	0.962	0.978
lsa:yes dim:300 sco:zeroone sw:yes	17	0.832	0.725	0.773	0.742	0.692	0.562
lsa:yes dim:300 sco:zeroone sw:-	18	0.995	0.989	0.986	0.981	0.966	0.971

Results - Evaluating Rankings

Spearman's Correlation

Blue: 1 (dark) to -1 (white)

Computed average and standard deviation of results, per model.

Spearman	LP1	LP2	LP3	LP4	LP5	LP6	LP7	LP8	LP9	LP10
lsa:- dim:- sco:count sw:yes	0.886	0.714	0.771	0.886	0.714	0.086	0.486	0.486	0.543	0.486
lsa:- dim:- sco:count sw:-	0.943	0.371	0.943	0.086	0.657	0.886	0.600	0.257	0.600	0.200
lsa:- dim:- sco:tfidf sw:yes	0.714	0.771	0.771	0.886	0.886	0.257	0.714	0.600	0.829	0.429
lsa:- dim:- sco:tfidf sw:-	0.771	0.657	0.771	0.029	0.657	0.943	0.714	0.257	0.714	0.543
lsa:- dim:- sco:zeroone sw:yes	0.714	0.714	0.543	0.886	0.771	0.257	0.771	0.486	0.257	0.429
lsa:- dim:- sco:zeroone sw:-	1.000	0.771	0.771	0.371	0.943	1.000	0.943	-0.086	0.829	0.714
lsa:yes dim:100 sco:count sw:yes	0.943	0.200	0.771	0.886	0.657	0.257	0.486	0.600	0.486	0.086
lsa:yes dim:100 sco:count sw:-	1.000	0.371	0.943	0.086	0.600	0.943	0.657	0.029	0.600	0.200
lsa:yes dim:100 sco:tfidf sw:yes	0.829	0.200	0.943	0.886	0.600	0.257	0.486	0.600	0.543	0.371
lsa:yes dim:100 sco:tfidf sw:-	1.000	0.600	0.943	0.029	0.657	0.943	0.600	0.371	0.714	0.200
lsa:yes dim:100 sco:zeroone sw:yes	0.829	0.314	0.829	0.829	0.829	0.600	0.543	0.257	0.600	0.257
lsa:yes dim:100 sco:zeroone sw:-	0.886	0.714	0.771	0.257	1.000	0.829	0.371	0.086	0.829	0.600
lsa:yes dim:300 sco:count sw:yes	0.771	0.200	0.886	0.886	0.657	0.257	0.257	0.657	0.486	0.486
lsa:yes dim:300 sco:count sw:-	1.000	0.371	0.943	0.086	0.600	0.943	0.600	0.029	0.600	0.200
lsa:yes dim:300 sco:tfidf sw:yes	0.943	0.314	0.943	0.829	0.657	0.600	0.257	0.657	0.486	0.486
lsa:yes dim:300 sco:tfidf sw:-	1.000	0.657	0.943	0.086	0.600	0.943	0.714	0.257	0.886	0.371
lsa:yes dim:300 sco:zeroone sw:yes	0.714	0.543	0.886	1.000	0.657	0.886	0.714	0.657	0.829	0.600
lsa:yes dim:300 sco:zeroone sw:-	0.943	0.886	0.714	0.371	1.000	0.943	0.657	-0.086	0.829	0.943

Results

I. Zero-one, sw, dim:300

II. TF-IDF, sw

III. TF-IDF, sw, dim:300

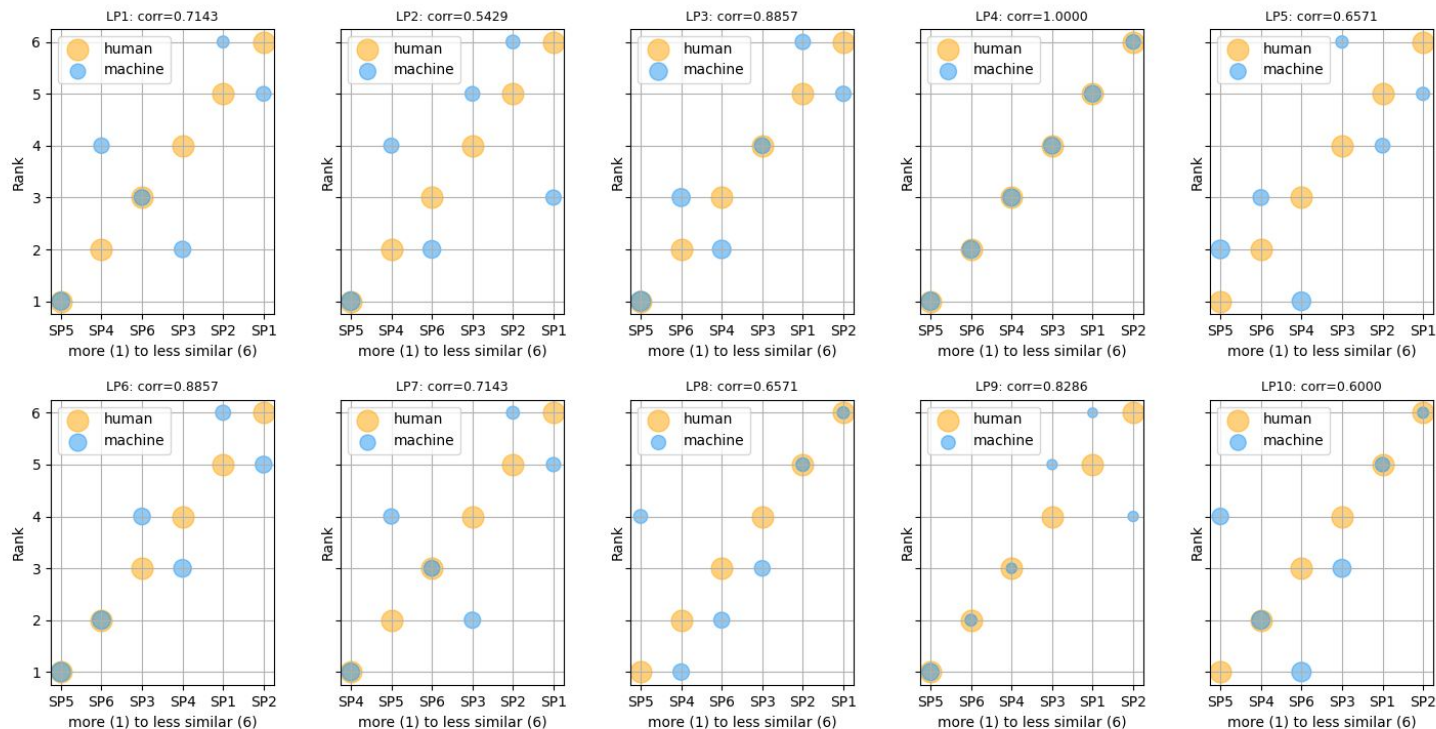
X. Zero-one, sw

Green: highest values on avg
Red: lowest values (0.3 up) on avg
Yellow: low values (0-0.3) on stdev

Model	Spearman			Kendall			O P Measure	
	AVG	ST DEV		AVG	ST DEV		AVG	ST DEV
Isa:- dim:- sco:count sw:yes	0.606	0.242		0.467	0.208		0.733	0.104
Isa:- dim:- sco:count sw:-	0.554	0.315		0.480	0.284		0.740	0.142
Isa:- dim:- sco:tfidf sw:yes (II)	0.686	0.204		0.533	0.201		0.767	0.101
Isa:- dim:- sco:tfidf sw:-	0.606	0.270		0.507	0.260		0.753	0.130
Isa:- dim:- sco:zeroone sw:yes (X)	0.583	0.222		0.427	0.189		0.713	0.095
Isa:- dim:- sco:zeroone sw:-	0.726	0.341		0.627	0.360		0.813	0.180
Isa:yes dim:100 sco:count sw:yes	0.537	0.291		0.413	0.296		0.707	0.148
Isa:yes dim:100 sco:count sw:-	0.543	0.360		0.480	0.358		0.740	0.179
Isa:yes dim:100 sco:tfidf sw:yes	0.571	0.256		0.427	0.267		0.713	0.134
Isa:yes dim:100 sco:tfidf sw:-	0.606	0.325		0.533	0.328		0.767	0.164
Isa:yes dim:100 sco:zeroone sw:yes	0.589	0.242		0.480	0.222		0.740	0.111
Isa:yes dim:100 sco:zeroone sw:-	0.634	0.300		0.520	0.275		0.760	0.138
Isa:yes dim:300 sco:count sw:yes	0.554	0.258		0.440	0.242		0.720	0.121
Isa:yes dim:300 sco:count sw:-	0.537	0.359		0.467	0.356		0.733	0.178
Isa:yes dim:300 sco:tfidf sw:yes (III)	0.617	0.240		0.507	0.244		0.753	0.122
Isa:yes dim:300 sco:tfidf sw:-	0.646	0.317		0.573	0.307		0.787	0.153
Isa:yes dim:300 sco:zeroone sw:yes (I)	0.749	0.146		0.613	0.183		0.807	0.091
Isa:yes dim:300 sco:zeroone sw:-	0.720	0.340		0.613	0.352		0.807	0.176

Results - Ranking Comparison. One model.

model-lsa-dim300-zeroone-stop



[Models on GitHub](#)

Graph of LP_x and rankings of SP(i,x), given by humans and the machine.

Circle area = similarity value.

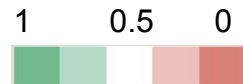
Conclusion and Further work

- Stopwords improve results. Each model has its own threshold
 - Best results: *LSA(dim:300), zero-one and stopwords*
1. **Human Ranking:** Larger population to obtain better human's perception
 2. **Test Data set:** Clear differentiation between similar sentences
 3. **Corpus:** Vocabulary and knowledge
 - Use a corpus that gives the system enough background on the topic to be tested
 4. **Generating our models:**
 - No syntactic information: "Computer Bag" or "Bag the computer", the same. Try n-grams
 5. **Storage:** Truncated LSA-models
 - Try Stemming and other preprocessing techniques to reduce matrix size



Thank you!

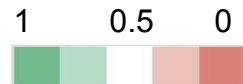
Results - Similarities - LP2



LP2	For the majority of women, breast cancer is not a death sentence. In the U.S., fewer women have been dying from breast cancer since the 1990s. Women whose breast cancer has not spread to other organs in their body like the lungs, liver, or brain, have over a 90% chance of being alive in 5 years. How long each person diagnosed with breast cancer will live and whether she will die from the disease depends on many things, including how much the cancer has spread throughout the body, the woman's overall health, and more.
SP(1,2)	Cancer are malignant cells that grow into more cells to form a tumor which needs to be treated to not become lethal.
SP(2,2)	Systemic therapy is for the whole body and can target and kill cancer in other parts of the body.
SP(3,2)	Depending on the advancement of breast cancer as well of other factors, the patient might survive.
SP(4,2)	The perception of breast cancer is changing everyday. Women are not afraid because they know that if the diagnosis is early, they will have the opportunity to cure it depending on the comorbidities.
SP(5,2)	For most of the women population who have been diagnosed with breast cancer, it is no longer a life threatening situation. Those whom their cancer has not spread to other organs like lungs, liver or brain have a 90% possibility of living the next five years or so.
SP(6,2)	Breast cancer mortality depends on many factors, like health, lifestyle, and spread to other organs, but cancer is less deadly than 20-30 years ago.

model		LP2					
		SP5	SP6	SP4	SP1	SP2	SP3
Isa:- dim:- sco:count sw:yes	1	0.754	0.749	0.579	0.502	0.449	0.299
Isa:- dim:- sco:count sw:-	2	0.995	0.993	0.991	0.982	0.995	0.986
Isa:- dim:- sco:tfidf sw:yes	3	0.678	0.664	0.524	0.393	0.401	0.366
Isa:- dim:- sco:tfidf sw:-	4	0.986	0.969	0.972	0.953	0.982	0.962
Isa:- dim:- sco:zeroone sw:yes	5	0.854	0.837	0.722	0.553	0.535	0.392
Isa:- dim:- sco:zeroone sw:-	6	0.991	0.977	0.975	0.968	0.973	0.968
Isa:yes dim:100 sco:count sw:yes	7	0.711	0.747	0.491	0.635	0.552	0.244
Isa:yes dim:100 sco:count sw:-	8	0.994	0.991	0.988	0.979	0.995	0.986
Isa:yes dim:100 sco:tfidf sw:yes	9	0.729	0.760	0.520	0.635	0.584	0.299
Isa:yes dim:100 sco:tfidf sw:-	10	0.988	0.979	0.976	0.963	0.987	0.970
Isa:yes dim:100 sco:zeroone sw:yes	11	0.899	0.897	0.768	0.846	0.795	0.720
Isa:yes dim:100 sco:zeroone sw:-	12	0.996	0.988	0.982	0.980	0.980	0.979
Isa:yes dim:300 sco:count sw:yes	13	0.681	0.714	0.483	0.608	0.562	0.288
Isa:yes dim:300 sco:count sw:-	14	0.994	0.991	0.988	0.978	0.995	0.985
Isa:yes dim:300 sco:tfidf sw:yes	15	0.673	0.631	0.452	0.580	0.548	0.332
Isa:yes dim:300 sco:tfidf sw:-	16	0.987	0.974	0.974	0.959	0.985	0.968
Isa:yes dim:300 sco:zeroone sw:yes	17	0.862	0.829	0.707	0.713	0.647	0.689
Isa:yes dim:300 sco:zeroone sw:-	18	0.993	0.982	0.979	0.973	0.977	0.973

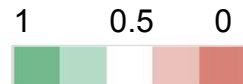
Results - Similarities - LP4



LP4	So far herbal products have not been shown to cure cancer. However, there are some herbs that may help patients deal with the side effects of cancer treatments. We recommend that you talk to your doctor before taking any vitamins and herbal products as they might have an effect on your cancer treatments.
SP(1,4)	Even though your attitude might not be a direct cause of cancer, it is helpful to stay positive and social.
SP(2,4)	There are many variables to determine whether a woman will die from breast cancer within the next five years.
SP(3,4)	Herbal products and vitamins are not linked to any treatment for cancer, however, it can minimize the secondary effects of proper treatment.
SP(4,4)	People are accustomed to consuming herbal products but there is no evidence of cure when they use them. However they could help with side effects of traditional cancer treatments.
SP(5,4)	So far, there is no proof that herbal products heal cancer, but there are some that can help with the treatment side effects. Nevertheless any herbal or vitamin products that you might want to try should be advised by your treating physician.
SP(6,4)	Herbal products may help with some cancer side effects, but you should discuss them with your doctor before starting to take any vitamins or supplements.

model		LP9					
		SP5	SP4	SP3	SP2	SP6	SP1
Isa:- dim:- sco:count sw:yes	1	0.694	0.520	0.477	0.460	0.409	0.388
Isa:- dim:- sco:count sw:-	2	0.996	0.996	0.996	0.984	0.997	0.994
Isa:- dim:- sco:tfidf sw:yes	3	0.548	0.385	0.356	0.301	0.413	0.394
Isa:- dim:- sco:tfidf sw:-	4	0.985	0.980	0.985	0.955	0.990	0.977
Isa:- dim:- sco:zeroone sw:yes	5	0.665	0.468	0.414	0.452	0.288	0.302
Isa:- dim:- sco:zeroone sw:-	6	0.982	0.973	0.973	0.968	0.986	0.962
Isa:yes dim:100 sco:count sw:yes	7	0.920	0.806	0.850	0.688	0.410	0.379
Isa:yes dim:100 sco:count sw:-	8	0.995	0.995	0.995	0.981	0.997	0.993
Isa:yes dim:100 sco:tfidf sw:yes	9	0.916	0.828	0.816	0.649	0.455	0.432
Isa:yes dim:100 sco:tfidf sw:-	10	0.988	0.986	0.989	0.963	0.993	0.981
Isa:yes dim:100 sco:zeroone sw:yes	11	0.940	0.749	0.868	0.760	0.843	0.678
Isa:yes dim:100 sco:zeroone sw:-	12	0.992	0.985	0.989	0.985	0.994	0.975
Isa:yes dim:300 sco:count sw:yes	13	0.891	0.707	0.833	0.650	0.441	0.393
Isa:yes dim:300 sco:count sw:-	14	0.995	0.995	0.995	0.981	0.997	0.993
Isa:yes dim:300 sco:tfidf sw:yes	15	0.822	0.667	0.734	0.604	0.478	0.436
Isa:yes dim:300 sco:tfidf sw:-	16	0.987	0.982	0.986	0.956	0.992	0.980
Isa:yes dim:300 sco:zeroone sw:yes	17	0.798	0.489	0.484	0.489	0.553	0.457
Isa:yes dim:300 sco:zeroone sw:-	18	0.986	0.976	0.979	0.973	0.990	0.965

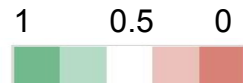
Results - Similarities - LP5



LP5	Breast cancer is a type of cancer that starts in the breast. Cancer starts when cells begin to grow out of control. It's important to understand that most breast lumps are benign and not cancer (malignant). Non-cancerous breast tumors are abnormal growths, but they do not spread outside of the breast. They are not life threatening, but some types of benign breast lumps can increase a woman's risk of getting breast cancer. Any breast lump or change needs to be checked by a healthcare professional to determine if it is benign or malignant (cancer) and if it might affect your future cancer risk.
SP(1,5)	Herbs might have beneficial effects on cancer treatments.
SP(2,5)	Getting good exercise, having a positive attitude, and letting your emotions out are good for those who have had cancer.
SP(3,5)	The growth of malignant cells on the breast is called breast cancer, however, not all are malignant but can indicate a predisposition.
SP(4,5)	Many types of breast tumor are benign. It's important to do a regular and frequent medical check up to figure out if it's malignant.
SP(5,5)	Cancer begins when body cells begin to grow out of control. This can happen in any part of the body, but when this happens in the breast tissue this is called breast cancer. Any breast lump or masses should be checked by a specialist to determine if it's malignant.
SP(6,5)	A lump or change in the breast may be cancer or it may not be -- those non cancerous lumps are not damaging to your health, but only a doctor can tell what kind of lump it is.

model		LP5					
		SP5	SP4	SP2	SP6	SP3	SP1
Isa:- dim:- sco:count sw:yes	1	0.900	0.897	0.555	0.534	0.517	0.279
Isa:- dim:- sco:count sw:-	2	0.995	0.997	0.990	0.997	0.996	0.945
Isa:- dim:- sco:tfidf sw:yes	3	0.675	0.663	0.448	0.482	0.371	0.271
Isa:- dim:- sco:tfidf sw:-	4	0.987	0.988	0.972	0.990	0.988	0.837
Isa:- dim:- sco:zeroone sw:yes	5	0.809	0.810	0.666	0.697	0.546	0.325
Isa:- dim:- sco:zeroone sw:-	6	0.993	0.990	0.982	0.993	0.990	0.938
Isa:yes dim:100 sco:count sw:yes	7	0.967	0.958	0.645	0.515	0.708	0.282
Isa:yes dim:100 sco:count sw:-	8	0.993	0.995	0.986	0.997	0.995	0.928
Isa:yes dim:100 sco:tfidf sw:yes	9	0.949	0.927	0.677	0.567	0.666	0.344
Isa:yes dim:100 sco:tfidf sw:-	10	0.988	0.991	0.977	0.993	0.991	0.847
Isa:yes dim:100 sco:zeroone sw:yes	11	0.955	0.918	0.785	0.847	0.892	0.701
Isa:yes dim:100 sco:zeroone sw:-	12	0.997	0.996	0.988	0.997	0.994	0.985
Isa:yes dim:300 sco:count sw:yes	13	0.950	0.941	0.564	0.504	0.651	0.297
Isa:yes dim:300 sco:count sw:-	14	0.993	0.995	0.986	0.997	0.995	0.930
Isa:yes dim:300 sco:tfidf sw:yes	15	0.904	0.865	0.583	0.552	0.634	0.386
Isa:yes dim:300 sco:tfidf sw:-	16	0.988	0.990	0.974	0.993	0.990	0.837
Isa:yes dim:300 sco:zeroone sw:yes	17	0.875	0.876	0.689	0.745	0.582	0.620
Isa:yes dim:300 sco:zeroone sw:-	18	0.995	0.992	0.984	0.995	0.991	0.947

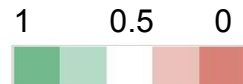
Results - Similarities - LP6



LP6	For most types of cancer, doctors use staging information to help plan treatment and to predict a person's outlook (prognosis). Although each person's situation is different, cancers with the same stage tend to have similar outlooks and are often treated the same way. The cancer stage is also a way for doctors to describe the extent of the cancer when they talk with each other about a person's cancer.
SP(1,6)	Even though not all new lumps or masses are malignant, they should be checked by a doctor to determine its risks.
SP(2,6)	Before taking herbal products, that may help some deal with cancer side effects, patients should consult with their doctor.
SP(3,6)	The advancement of cancer is categorized and described in stages, which also works as a guide for treatment and prediction.
SP(4,6)	Doctors usually use stages for cancer so they can predict the evolution and can choose the treatment. Actually they can describe the extension of the tumor through staging.
SP(5,6)	When doctors are dealing with cancer patients, Oncologists use staging information as a reference to determine the severity of the patient, and their treatment as well. Even though each patient has their own characteristics, there are similar symptoms during each disease's stage.
SP(6,6)	Each person with cancer is unique, but cancer stages allow doctors to group similar patients together when they discuss treatment options and outcomes.

model		LP6					
		SP5	SP4	SP2	SP1	SP6	SP3
Isa:- dim:- sco:count sw:yes	1	0.888	0.615	0.602	0.579	0.526	0.384
Isa:- dim:- sco:count sw:-	2	0.996	0.986	0.966	0.986	0.994	0.994
Isa:- dim:- sco:tfidf sw:yes	3	0.720	0.499	0.463	0.356	0.462	0.335
Isa:- dim:- sco:tfidf sw:-	4	0.987	0.969	0.905	0.962	0.977	0.985
Isa:- dim:- sco:zeroone sw:yes	5	0.861	0.716	0.709	0.645	0.706	0.497
Isa:- dim:- sco:zeroone sw:-	6	0.988	0.979	0.947	0.977	0.981	0.979
Isa:yes dim:100 sco:count sw:yes	7	0.945	0.690	0.722	0.569	0.691	0.521
Isa:yes dim:100 sco:count sw:-	8	0.996	0.986	0.959	0.984	0.992	0.993
Isa:yes dim:100 sco:tfidf sw:yes	9	0.900	0.768	0.775	0.535	0.739	0.569
Isa:yes dim:100 sco:tfidf sw:-	10	0.992	0.975	0.944	0.969	0.987	0.988
Isa:yes dim:100 sco:zeroone sw:yes	11	0.939	0.886	0.896	0.752	0.928	0.880
Isa:yes dim:100 sco:zeroone sw:-	12	0.996	0.988	0.969	0.986	0.993	0.985
Isa:yes dim:300 sco:count sw:yes	13	0.907	0.620	0.647	0.548	0.644	0.479
Isa:yes dim:300 sco:count sw:-	14	0.996	0.986	0.958	0.984	0.992	0.993
Isa:yes dim:300 sco:tfidf sw:yes	15	0.803	0.655	0.664	0.435	0.689	0.542
Isa:yes dim:300 sco:tfidf sw:-	16	0.990	0.973	0.914	0.963	0.984	0.987
Isa:yes dim:300 sco:zeroone sw:yes	17	0.900	0.831	0.785	0.710	0.845	0.785
Isa:yes dim:300 sco:zeroone sw:-	18	0.991	0.985	0.956	0.979	0.988	0.982

Results - Similarities - LP9



LP9	A mastectomy is a type of surgery to remove the entire breast. There are several types of mastectomies. Some mastectomies remove all of the breast tissue, including the nipple and areola and the other types of mastectomies also remove some of the nearby lymph nodes
SP(1,9)	A lumpectomy is surgery to remove the cancer tumor from the breast.
SP(2,9)	It is difficult to determine by a new lumps shape or hardness whether it is cancer or not.
SP(3,9)	There are different types of mastectomy, depending on how much of the breast they remove, up to the complete removal.
SP(4,9)	Mastectomy is a way to get rid of breast cancer and it would be as radical depending on how far the cancer has spread. The surgeons may take out tissue as long is compromised with the disease.
SP(5,9)	The mastectomy is a surgery that removes all the breast tissue. There are different types of mastectomy depending on the area removed. The first one focuses on all breast tissues including the nipple and areola, while others also add nearby lymph nodes to the list.
SP(6,9)	Mastectomies are surgeries that remove some or all of the breast, and may also remove adjacent lymph nodes.

model		LP9					
		SP5	SP4	SP3	SP2	SP6	SP1
Isa:- dim:- sco:count sw:yes	1	0.694	0.520	0.477	0.460	0.409	0.388
Isa:- dim:- sco:count sw:-	2	0.996	0.996	0.996	0.984	0.997	0.994
Isa:- dim:- sco:tfidf sw:yes	3	0.548	0.385	0.356	0.301	0.413	0.394
Isa:- dim:- sco:tfidf sw:-	4	0.985	0.980	0.985	0.955	0.990	0.977
Isa:- dim:- sco:zeroone sw:yes	5	0.665	0.468	0.414	0.452	0.288	0.302
Isa:- dim:- sco:zeroone sw:-	6	0.982	0.973	0.973	0.968	0.986	0.962
Isa:yes dim:100 sco:count sw:yes	7	0.920	0.806	0.850	0.688	0.410	0.379
Isa:yes dim:100 sco:count sw:-	8	0.995	0.995	0.995	0.981	0.997	0.993
Isa:yes dim:100 sco:tfidf sw:yes	9	0.916	0.828	0.816	0.649	0.455	0.432
Isa:yes dim:100 sco:tfidf sw:-	10	0.988	0.986	0.989	0.963	0.993	0.981
Isa:yes dim:100 sco:zeroone sw:yes	11	0.940	0.749	0.868	0.760	0.843	0.678
Isa:yes dim:100 sco:zeroone sw:-	12	0.992	0.985	0.989	0.985	0.994	0.975
Isa:yes dim:300 sco:count sw:yes	13	0.891	0.707	0.833	0.650	0.441	0.393
Isa:yes dim:300 sco:count sw:-	14	0.995	0.995	0.995	0.981	0.997	0.993
Isa:yes dim:300 sco:tfidf sw:yes	15	0.822	0.667	0.734	0.604	0.478	0.436
Isa:yes dim:300 sco:tfidf sw:-	16	0.987	0.982	0.986	0.956	0.992	0.980
Isa:yes dim:300 sco:zeroone sw:yes	17	0.798	0.489	0.484	0.489	0.553	0.457
Isa:yes dim:300 sco:zeroone sw:-	18	0.986	0.976	0.979	0.973	0.990	0.965

Results - Kendall Tau Correlation

Kendall Tau	LP1	LP2	LP3	LP4	LP5	LP6	LP7	LP8	LP9	LP10
Isa:- dim:- sco:count sw:yes	0.733	0.467	0.600	0.733	0.600	0.067	0.333	0.333	0.467	0.333
Isa:- dim:- sco:count sw:-	0.867	0.333	0.867	0.067	0.600	0.733	0.467	0.200	0.467	0.200
Isa:- dim:- sco:tfidf sw:yes	0.467	0.600	0.600	0.733	0.733	0.200	0.600	0.467	0.733	0.200
Isa:- dim:- sco:tfidf sw:-	0.600	0.600	0.600	-0.067	0.600	0.867	0.600	0.200	0.600	0.467
Isa:- dim:- sco:zeroone sw:yes	0.467	0.467	0.467	0.733	0.600	0.200	0.600	0.333	0.200	0.200
Isa:- dim:- sco:zeroone sw:-	1.000	0.600	0.600	0.333	0.867	1.000	0.867	-0.200	0.600	0.600
Isa:yes dim:100 sco:count sw:yes	0.867	0.067	0.600	0.733	0.600	0.200	0.333	0.467	0.333	-0.067
Isa:yes dim:100 sco:count sw:-	1.000	0.333	0.867	0.067	0.467	0.867	0.600	-0.067	0.467	0.200
Isa:yes dim:100 sco:tfidf sw:yes	0.600	0.067	0.867	0.733	0.467	0.200	0.333	0.467	0.467	0.067
Isa:yes dim:100 sco:tfidf sw:-	1.000	0.467	0.867	-0.067	0.600	0.867	0.467	0.333	0.600	0.200
Isa:yes dim:100 sco:zeroone sw:yes	0.600	0.200	0.733	0.733	0.733	0.467	0.467	0.200	0.467	0.200
Isa:yes dim:100 sco:zeroone sw:-	0.733	0.467	0.600	0.200	1.000	0.733	0.333	0.067	0.600	0.467
Isa:yes dim:300 sco:count sw:yes	0.600	0.067	0.733	0.733	0.600	0.200	0.200	0.600	0.333	0.333
Isa:yes dim:300 sco:count sw:-	1.000	0.333	0.867	0.067	0.467	0.867	0.467	-0.067	0.467	0.200
Isa:yes dim:300 sco:tfidf sw:yes	0.867	0.200	0.867	0.600	0.600	0.467	0.200	0.600	0.333	0.333
Isa:yes dim:300 sco:tfidf sw:-	1.000	0.600	0.867	0.067	0.467	0.867	0.600	0.200	0.733	0.333
Isa:yes dim:300 sco:zeroone sw:yes	0.467	0.467	0.733	1.000	0.467	0.733	0.467	0.600	0.733	0.467
Isa:yes dim:300 sco:zeroone sw:-	0.867	0.733	0.467	0.333	1.000	0.867	0.600	-0.200	0.600	0.867

Results - Order Preservation Measure

Order Preservation Measure	LP1	LP2	LP3	LP4	LP5	LP6	LP7	LP8	LP9	LP10
Isa:- dim:- sco:count sw:yes	0.867	0.733	0.800	0.867	0.800	0.533	0.667	0.667	0.733	0.667
Isa:- dim:- sco:count sw:-	0.933	0.667	0.933	0.533	0.800	0.867	0.733	0.600	0.733	0.600
Isa:- dim:- sco:tfidf sw:yes	0.733	0.800	0.800	0.867	0.867	0.600	0.800	0.733	0.867	0.600
Isa:- dim:- sco:tfidf sw:-	0.800	0.800	0.800	0.467	0.800	0.933	0.800	0.600	0.800	0.733
Isa:- dim:- sco:zeroone sw:yes	0.733	0.733	0.733	0.867	0.800	0.600	0.800	0.667	0.600	0.600
Isa:- dim:- sco:zeroone sw:-	1.000	0.800	0.800	0.667	0.933	1.000	0.933	0.400	0.800	0.800
Isa:yes dim:100 sco:count sw:yes	0.933	0.533	0.800	0.867	0.800	0.600	0.667	0.733	0.667	0.467
Isa:yes dim:100 sco:count sw:-	1.000	0.667	0.933	0.533	0.733	0.933	0.800	0.467	0.733	0.600
Isa:yes dim:100 sco:tfidf sw:yes	0.800	0.533	0.933	0.867	0.733	0.600	0.667	0.733	0.733	0.533
Isa:yes dim:100 sco:tfidf sw:-	1.000	0.733	0.933	0.467	0.800	0.933	0.733	0.667	0.800	0.600
Isa:yes dim:100 sco:zeroone sw:yes	0.800	0.600	0.867	0.867	0.867	0.733	0.733	0.600	0.733	0.600
Isa:yes dim:100 sco:zeroone sw:-	0.867	0.733	0.800	0.600	1.000	0.867	0.667	0.533	0.800	0.733
Isa:yes dim:300 sco:count sw:yes	0.800	0.533	0.867	0.867	0.800	0.600	0.600	0.800	0.667	0.667
Isa:yes dim:300 sco:count sw:-	1.000	0.667	0.933	0.533	0.733	0.933	0.733	0.467	0.733	0.600
Isa:yes dim:300 sco:tfidf sw:yes	0.933	0.600	0.933	0.800	0.800	0.733	0.600	0.800	0.667	0.667
Isa:yes dim:300 sco:tfidf sw:-	1.000	0.800	0.933	0.533	0.733	0.933	0.800	0.600	0.867	0.667
Isa:yes dim:300 sco:zeroone sw:yes	0.733	0.733	0.867	1.000	0.733	0.867	0.733	0.800	0.867	0.733
Isa:yes dim:300 sco:zeroone sw:-	0.933	0.867	0.733	0.667	1.000	0.933	0.800	0.400	0.800	0.933